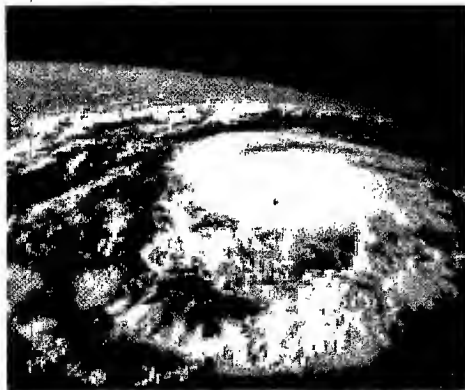
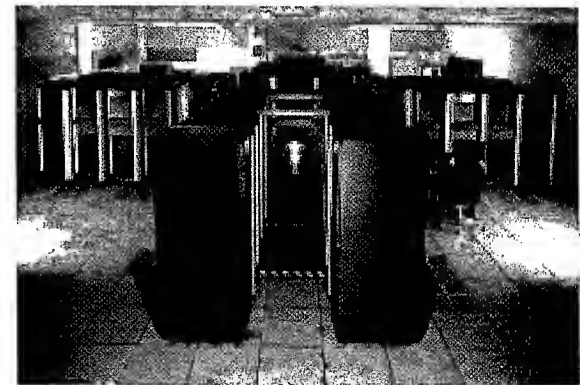


## *Scaling NASA Applications to 1024 CPUs on Origin 3K*



*NRC Review*  
*June 13, 2002*

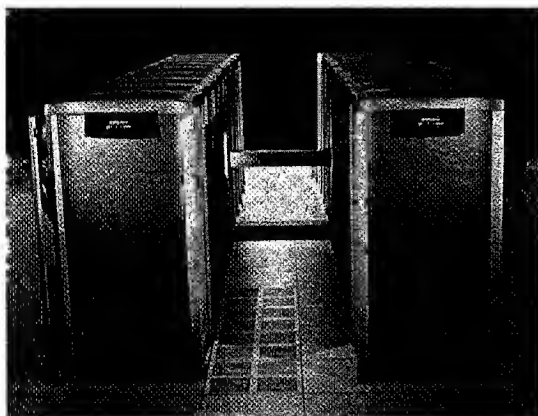
*Jim Taft*  
*NASA AMES Research Center*  
*[jtaft@nas.nasa.gov](mailto:jtaft@nas.nasa.gov)*



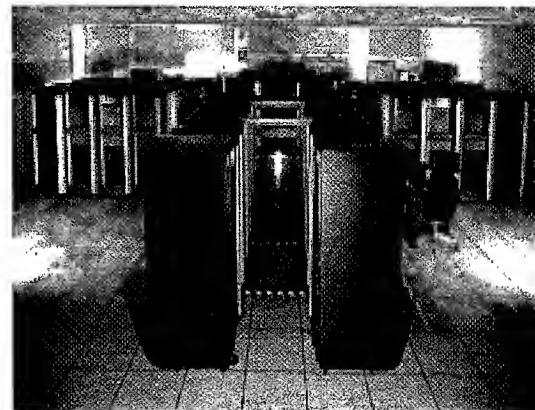
# ***The Origin SSI Revolution at NASA Ames (1998-2001)***



***Steger: O2K 256CPUs 250MHz***



***Lomax: O2K 512CPUs 400MHz***



***Chapman: O3K 1024CPUs 600MHz***





## ***MLP a History of Success***



***The long and highly successful joint SGI-NASA research effort in ever larger SSI systems was to a large degree the result of the successful development of the MLP scalable parallel programming paradigm developed at ARC:***

- MLP scaling in real production codes justified ever larger systems at NAS***
- MLP scaling on 256p Origin 2000 gave SGI impetus to productize 256p***
- MLP scaling on 512 gave SGI courage to build 1024p O3K***
- History of MLP success resulted in IBM Star Cluster based MLP effort***

# ***MLPlib - Simple***



## ***Scalable Parallelism***

### ***at NASA Ames***



## ***Performance - The Focus is on Parallelism***



***Parallelism is the key to performance on any system manufactured today. If you don't scale to hundreds of CPUs, you won't get to the 100+ GFLOPS you need today to stay competitive in high end computing. Parallelism was being aggressively pursued on two fronts. Now there are three.***

- ***Message Passing Interface (MPI)***
  - ***Arcane and complex user interface - 100 routines, 50,000+ lines of source***
  - ***Explicit "messages" – large latencies – very slow***
  - ***User provides all parallel decomposition/code modification***
  - ***Often requires simplification of physics for scaling***
- ***Shared Memory Parallelism (OpenMP)***
  - ***Really acceptable only for small processor counts***
  - ***Very difficult to scale to 100's of CPU's without major rewrite***
- ***NASA's Shared Memory Multi-Level Parallelism (MLP)***
  - ***Simple extension to Cray parallel/vector programming model - 3 routines, 150 lines of source***
  - ***No messaging - All communication via shared memory***
  - ***Much easier to build/port code than MPI (Man months vs. Man years)***
  - ***Minimum changes OVERFLOW (MPI/MLP=20,000/800 lines), FVCORE (8000/400 lines)***
  - ***Dramatically better performance with increasing processor count***

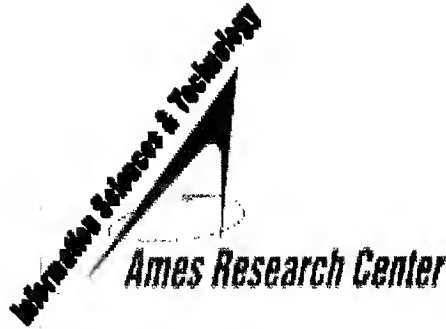


## ***What is MLP?***



***Shared Memory Multi-level Parallelism (MLP) is the utilization of multiple levels of parallelism within an application executing on a NUMA based system architecture in order to increase its parallel efficiency during execution. It is an open system design (runs on any SMP) and has the following attributes:***

- ***Two levels of parallelism (the so-called “hybrid” approach)***
- ***Coarse grained parallelism provided by Unix forked processes***
- ***Fine grained parallelism provided by the compiler at loop level (OpenMP)***
- ***No messaging - communication through “global” common blocks***
- ***Targeted for the new large CPU count NUMA SMP systems***
- ***But method has been adapted to execute across clusters as well***



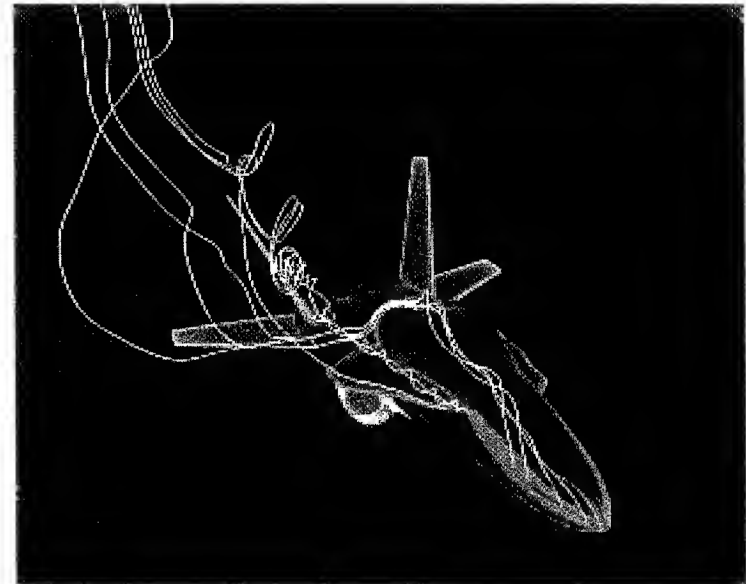
## *How do you use it - MLPlib*



*MLPlib consists of 3 routines - how hard can it be?*

- ***Subroutine MLP\_GETMEM(numvar,ipoint,ysize)***
  - *Called once only from main program*
  - *Allocates the global common blocks*
- ***Subroutine MLP\_FORKIT(numpro,nowpro,numcps,idopin)***
  - *Called once only from main program*
  - *Forks the processes*
- ***Subroutine MLP\_BARRIER(nowpro,numpro)***
  - *Called as often as needed*
  - *Barrier synchs the processes*

***CFD at NASA***



***OVERFLOW-MLP***





## **MLP - A New Concept for Multi-zonal CFD**



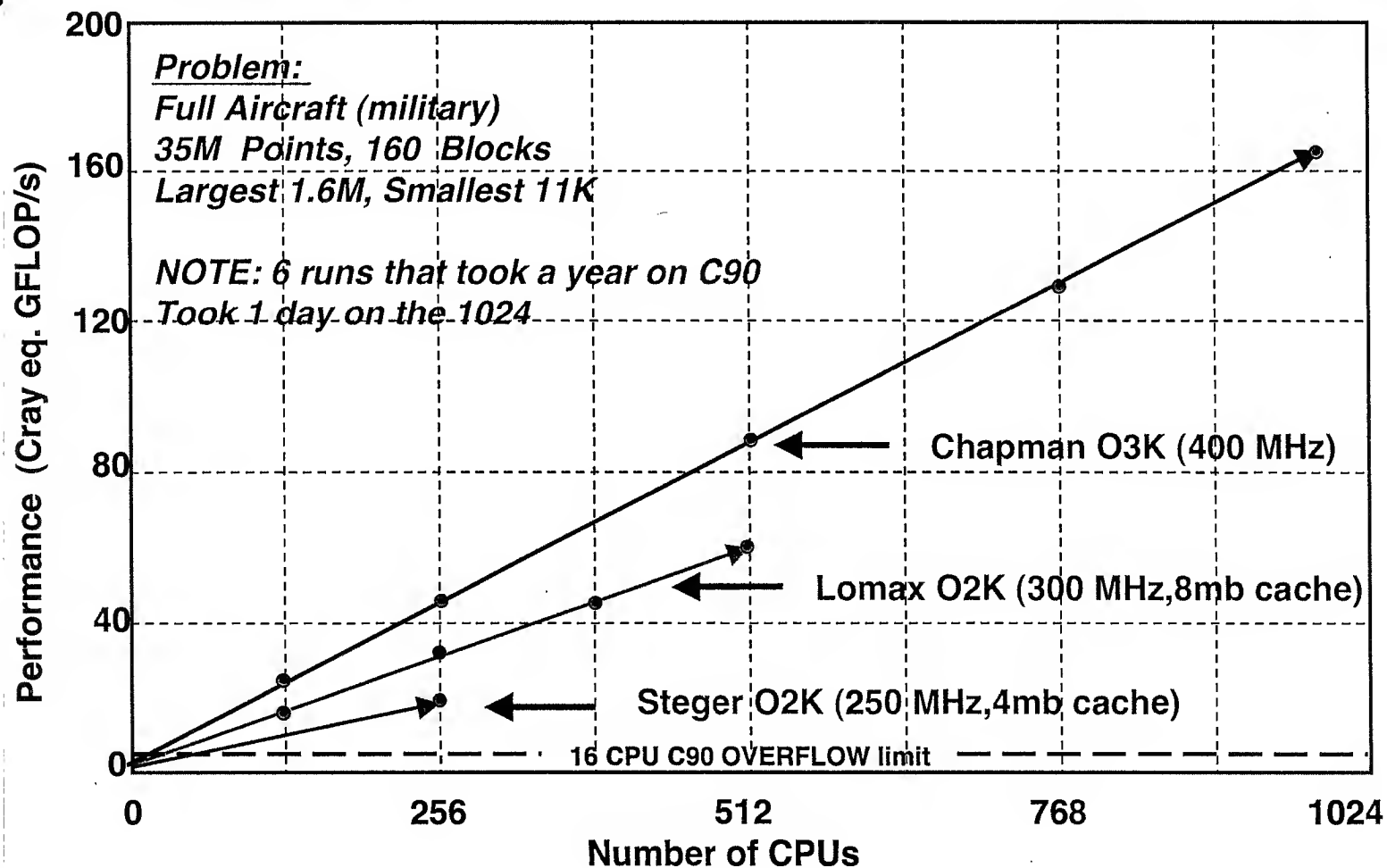
***NASA's Multi-zonal CFD codes like OVERFLOW, CFL3D, LAURA, INS3D, and TLNS3D to name a few, are ideal candidates for MLP parallelism. These codes decompose a large region of interest into many linked smaller 3D regions. These smaller regions can be solved mostly in parallel, with the occasional exchange of boundary information at the end of a time step.***

***In short, the recipe for converting a multi-zonal CFD code to MLP is:***

- Spawn MLP parallel processes***
- Assign groups of 3D zones to each MLP process***
- Solve the groups of zones in parallel***
- Assign groups of CPUs to each MLP process***
- Use the CPUs in a group for fine grained parallelism for each zone***
- Use shared memory arenas to hold all global data (BCs etc)***
- Synchronize computation as needed with barriers***

# OVERFLOW-MLP Performance vs CPU Count

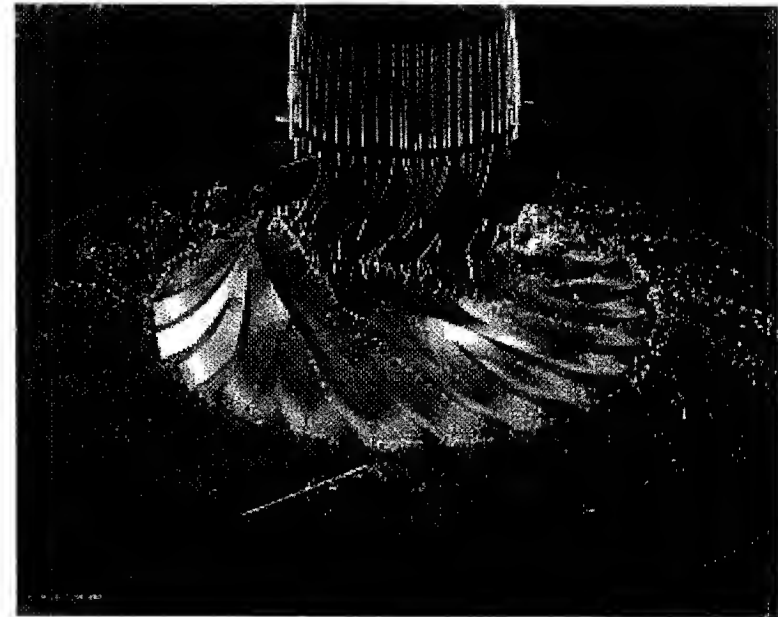
Systems: 1024 CPU O3K , 256&512 CPU O2KS



***Turbomachinery***

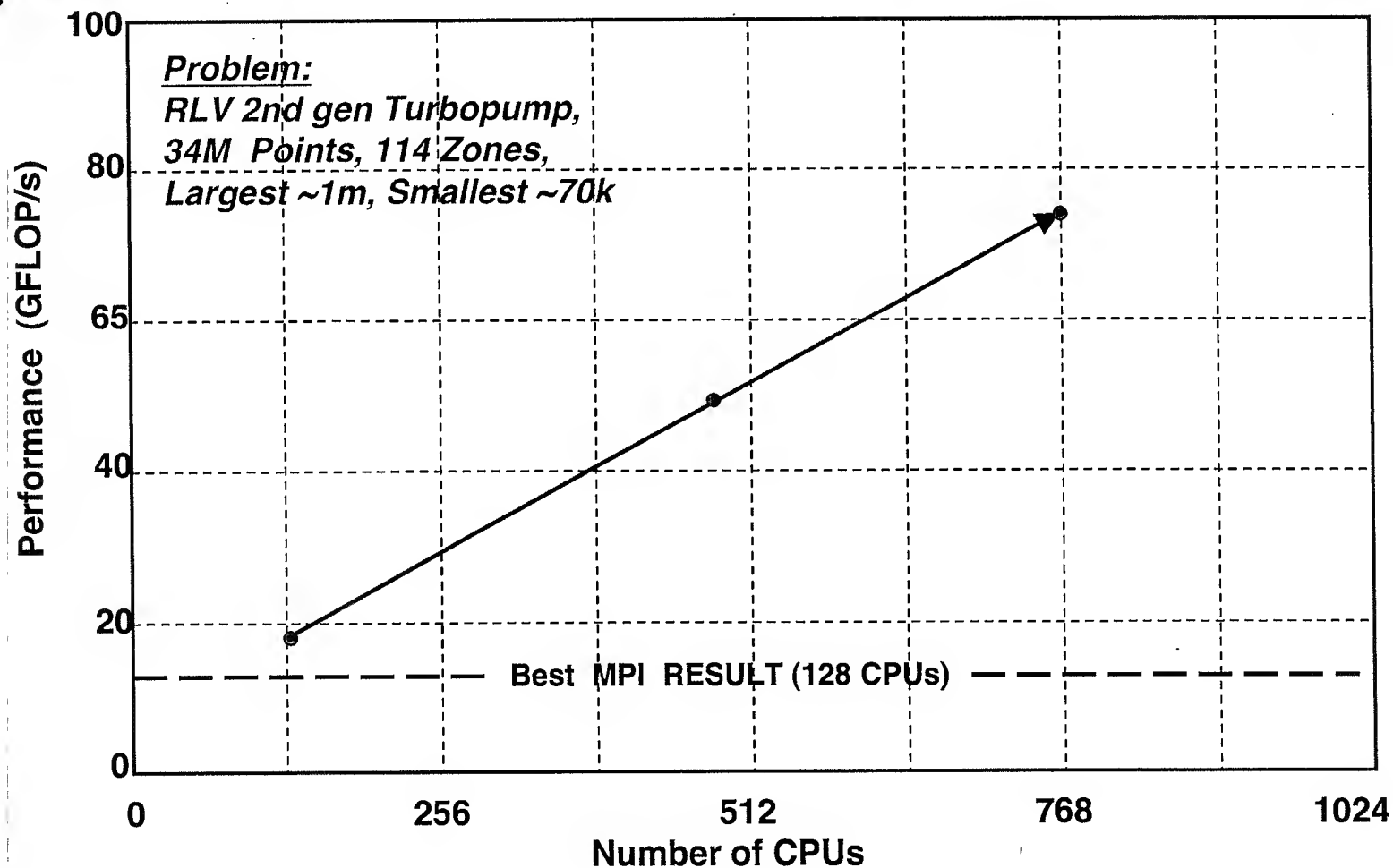
***CFD at NASA***

***INS3D-MLP***

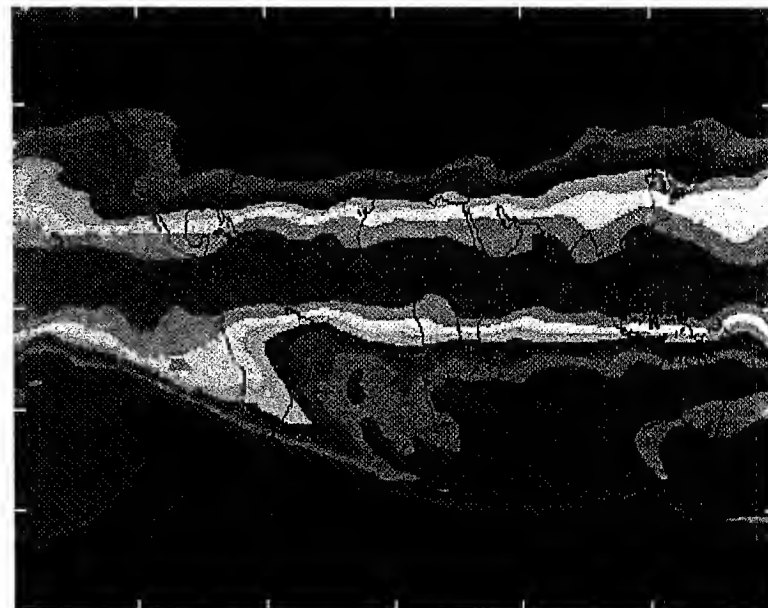


# INS3D-MLP Performance vs CPU Count

System: 1024 CPU O3K (400 MHz)



***Climate Modeling***  
***at NASA Ames***



***GEOS4-MLP***



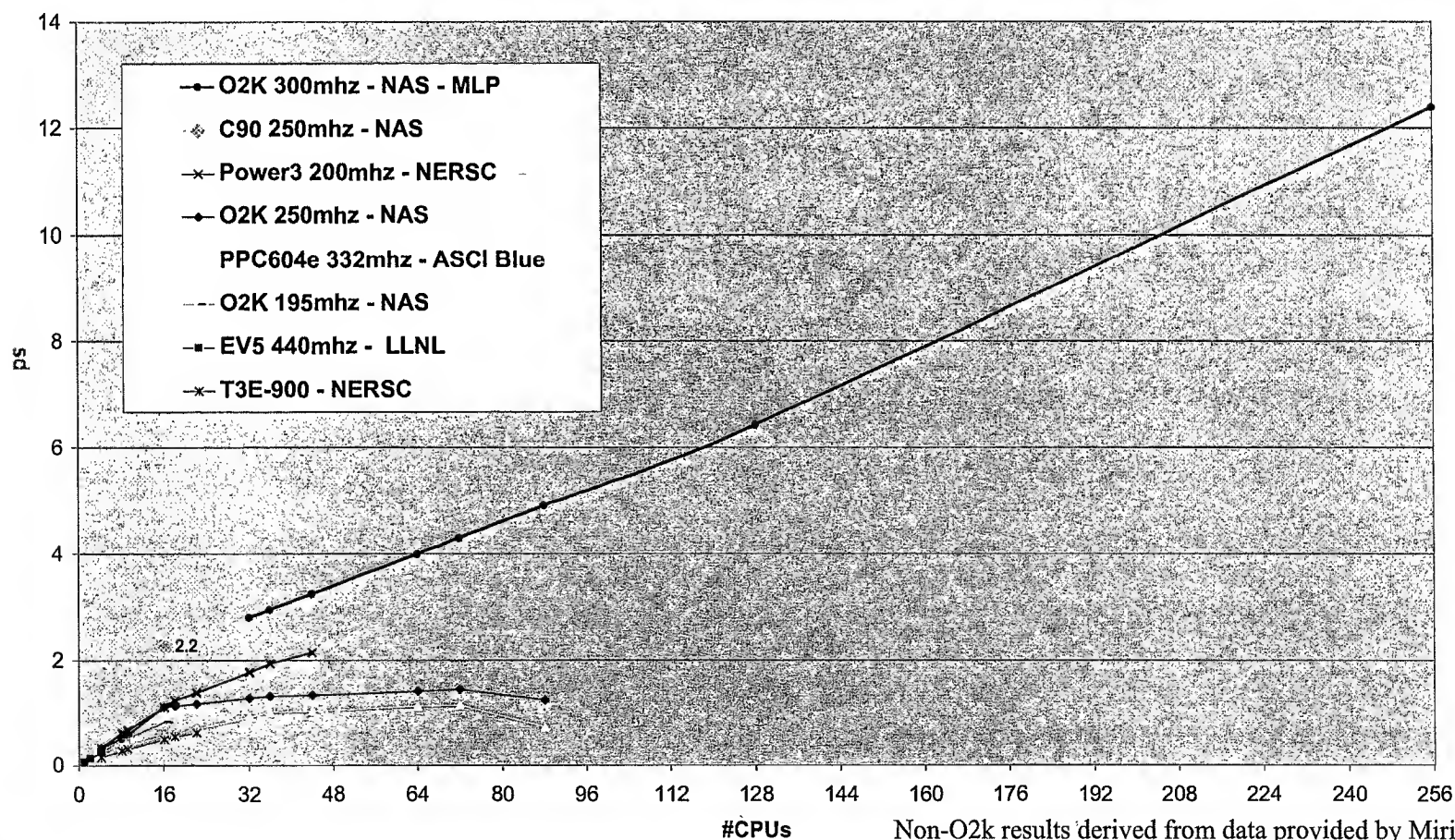
# NASA Climate Modeling

## Initial FVCORE-MLP Scaling on Popular Systems 06/00

(NAS MLP work completed in 21days)



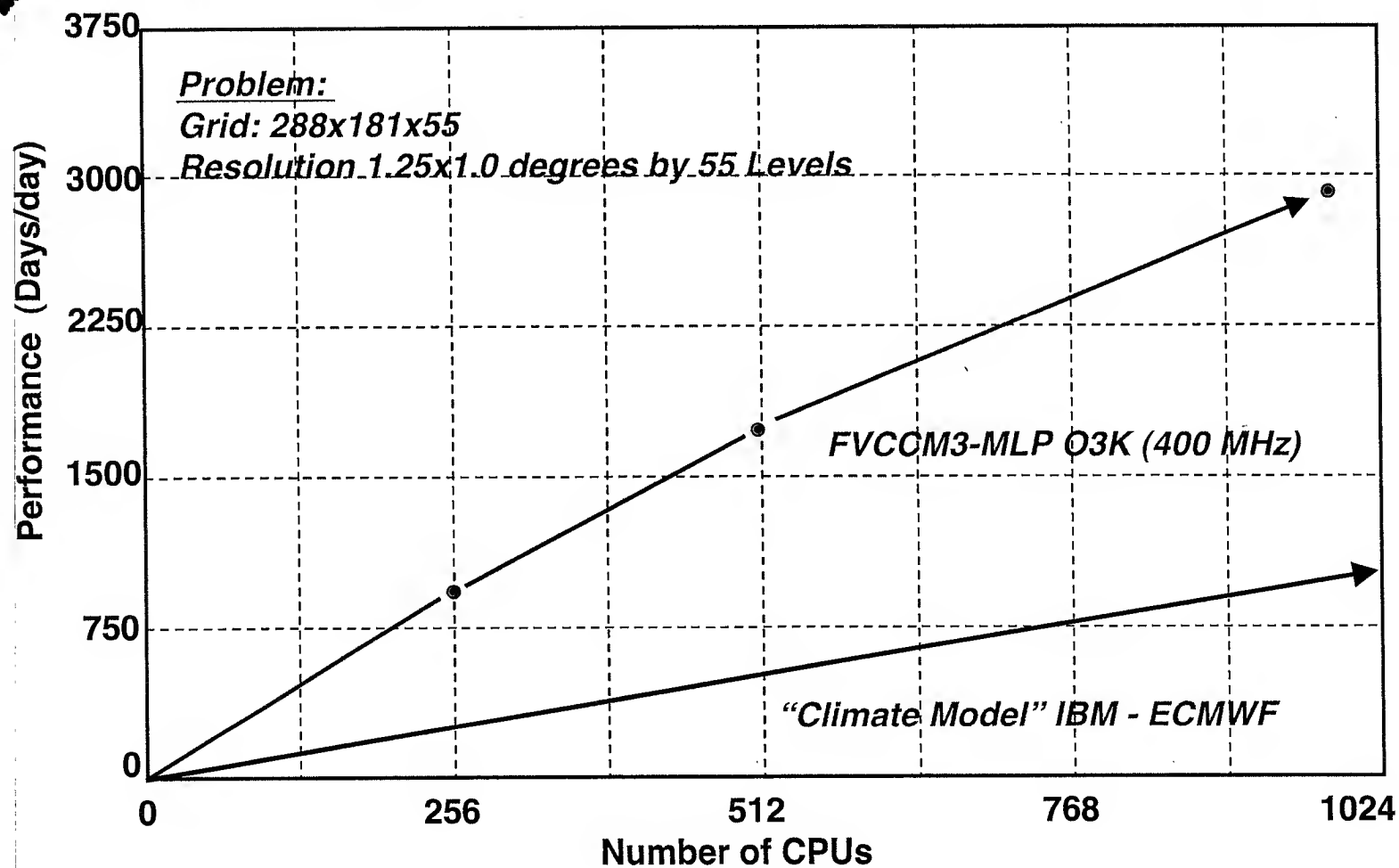
FVCORE 2x2.5@55 Levels (B55)



Non-O2k results derived from data provided by Mirin et al

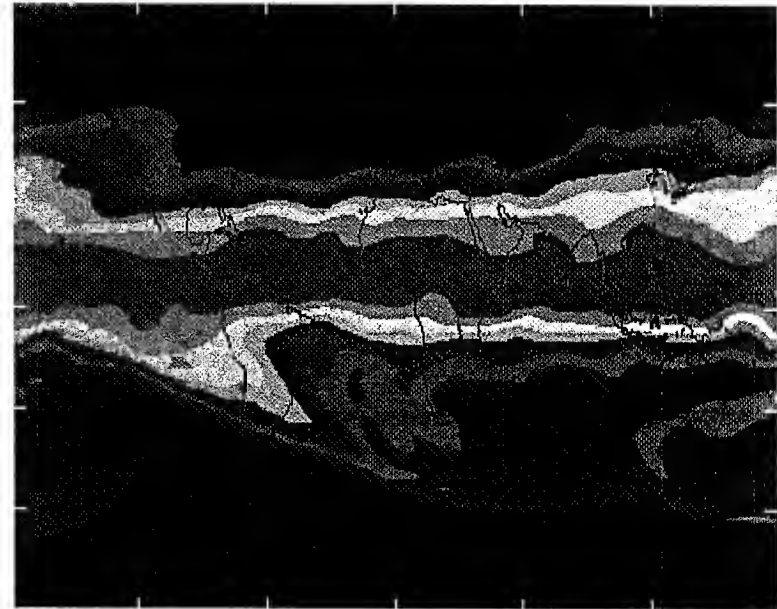
## FVCCM3-MLP Performance vs CPU Count (2/02)

System: 1024 CPU O3K 400 MHz (chapman)



# ***Parallel Ocean***

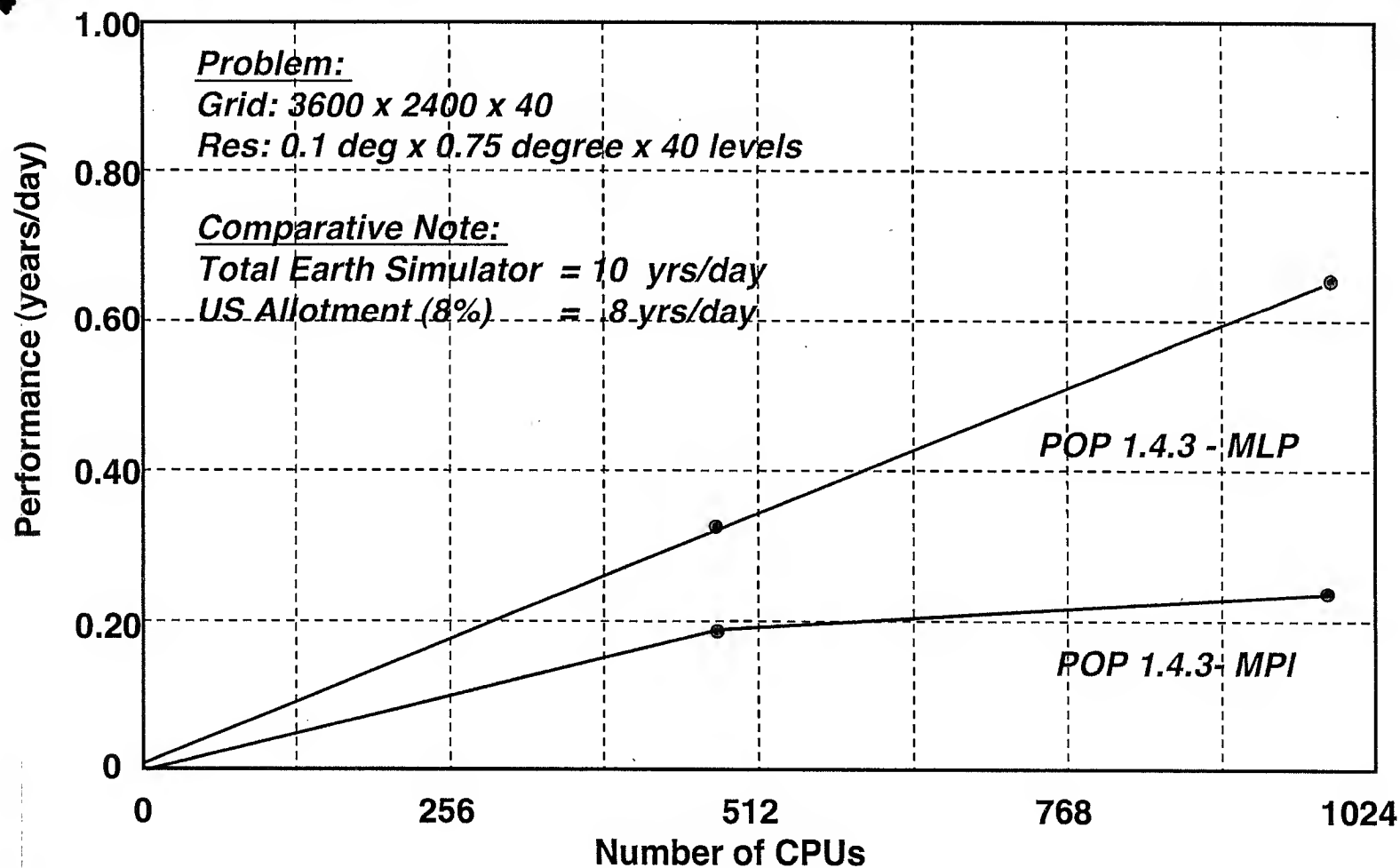
## ***Modeling - POP***





# POP Performance vs CPU Count (6/02)

System: 1024 CPU O3K 600 MHz (chapman)





## ***The Future of HPC at NASA***



***At a November, 2001 presentation to Dan Goldin we were asked:  
Why haven't you ordered an even larger system?***

***He asks this question because:***

***The scalable software techniques are in place  
MLP parallelism at NASA Ames scales for real science - not kernels  
MLP is a practical means for migrating codes in our lifetime - weeks***

***The hardware architecture exists  
SSI NUMA architectures let scientists do science - NOT MPI DEBUG  
Migration to larger SSI systems is transparent to users  
MLP+NUMA supports "Cray programming model", the "standard"***

***Cost is more than competitive  
SSI NUMA supports higher fractions of sustained vs peak.  
Time to market of applications is years earlier***



## ***So what got Accelerated by NASA Ames?***



### ***Hardware:***

***SGI developed 256, 512, and 1024 CPU systems not in their plan***

### ***MLPlib:***

***A new very fast way of getting to large parallel scaling unlike MPI***

***Maintains shared memory "Cray programming model" for users***

### ***True Production codes and their Accelerated Performance:***

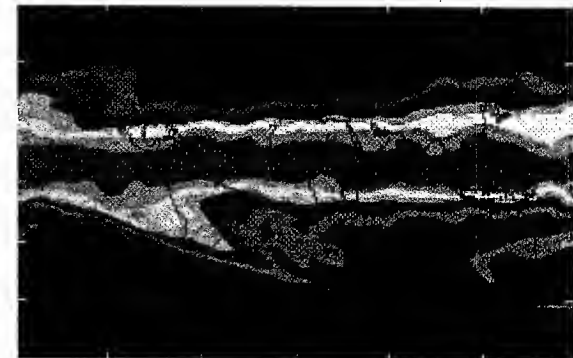
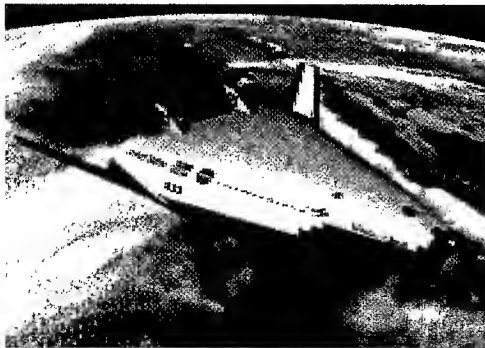
<u><b>Code</b></u>	<u><b>Speedup</b></u>	<u><b>Effort</b></u>	<u><b>Code size</b></u>
<b>OVERFLOW (CFD)</b>	<b>35X</b>	<b>1 man-year</b>	<b>100,000 lines</b>
<b>INS3D (CFD)</b>	<b>33X</b>	<b>2 man-years</b>	<b>50,000 lines</b>
<b>COSMOS core (MD)</b>	<b>35x</b>	<b>6 man-months</b>	<b>30,000 lines</b>
<b>DAS (Climate)</b>	<b>10x</b>	<b>2.5 man-years</b>	<b>350,000 lines</b>
<b>FVCCM3 (Climate)</b>	<b>40x</b>	<b>4 man-months</b>	<b>65,000 lines</b>

***Note: MPI versions of codes above took about 5-10x longer to build***

***THIS IS THE LEGACY OF MLP SHARED MEMORY PROGRAMMING***

# ***The Next Step in MLP Performance on SGI's Origin 1024p SSI System***

## ***Climate and CFD***





## ***Optimization Efforts - CFD/Climate Codes***



The bulk of NASA's NAS supercomputing cycles remain CFD related, with additional cycles devoted to climate and molecular modeling. This year, the NAS-ACL will focus on resolving remaining barriers to higher performance in CFD/Climate on the SGI shared memory Origin platforms. Major efforts will be:

- o **OVERFLOW-MLP - Sustain 250 GFLOP/s on 1024p system**
  - o Increase single CPU performance
  - o Reduce overall memory size and traffic
  - o Important because almost interactive electronic wind tunnel
- o **FVCCM3-MLP - Sustain 8000 days/day on 1024p system**
  - o Reduce communication burden
  - o Increase single CPU performance
  - o Important because same as Fujitsu 100 CPU result (ECMWF)

## Star Cluster - The Next Step in MLP

